

Non-Bayesian Additive Regularization for Multimodal Topic Modeling of Large Collections

Konstantin Vorontsov
Yandex, Moscow Institute
of Physics and Technology
voron@yandex-team.ru

Oleksandr Frei
Schlumberger
Information Solutions
oleksandr.frei@gmail.com

Murat Apishev
Moscow State University
great-mel@yandex.ru

Peter Romov
Yandex
peter@romov.ru

Marina Dudarenko
Moscow State University
m.dudarenko@gmail.com

ABSTRACT

Probabilistic topic modeling of text collections is a powerful tool for statistical text analysis based on the preferential use of graphical models and Bayesian learning. Additive regularization for topic modeling (ARTM) is a recent semi-probabilistic approach, which provides a much simpler inference for many models previously studied only in the Bayesian settings. Additive regularization makes topic models easier to design, infer, combine, and explain, thus reducing barriers to entry into topic modeling research field. In this paper we develop the ARTM approach in three directions: multimodal data modeling, online algorithm, and parallel implementation. We announce the BigARTM open source project (<http://bigartm.org>) for regularized multimodal topic modeling of large collections. Experiments on Wikipedia corpus show that BigARTM performs faster and gives better perplexity comparing to other popular packages, such as Vowpal Wabbit and Gensim. We also demonstrate several unique BigARTM features, such as additive combination of regularizers, topic sparsing and decorrelation, multimodal and multilanguage modeling, which are not available in the other software packages for topic modeling.

Keywords

Probabilistic Topic Modeling, Probabilistic Latent Semantic Analysis, Latent Dirichlet Allocation, Additive Regularization of Topic Models, Stochastic Matrix Factorization, EM-algorithm, BigARTM.

1. INTRODUCTION

Topic modeling is a rapidly developing branch of statistical text analysis [2]. Topic model reveals a hidden thematic structure of a text collection and finds a compressed representation of each document in terms of its topics. Practical applications of topic models include many areas, such as in-

formation retrieval for long-text queries, classification, categorization, summarization and segmentation of texts. Topic models are increasingly used for non-textual and heterogeneous data including signals, images, video and networks. More ideas, models and applications are outlined in the survey [7].

From a statistical point of view, a probabilistic topic model (PTM) defines each topic by a multinomial distribution over words, and then describes each document with a multinomial distribution over topics.

From an optimizational point of view, topic modeling can be considered as a special case of approximate stochastic matrix factorization. To learn a factorized representation of a text collection is an ill-posed problem, which has an infinite set of solutions. A typical regularization approach in this case is to impose problem-specific constraints in a form of additive penalty terms in the optimization criterion.

Modern literature on topic modeling offers hundreds of models adapted to different situations. Nevertheless, most of these models are too difficult for practitioners to quickly understand, adapt and embed into applications. This leads to a common practice of tasting only the basic out-of-date models such as *Probabilistic Latent Semantic Analysis*, PLSA [12] and *Latent Dirichlet Allocation*, LDA [4]. Most practical inconveniences are rooted in Bayesian learning, which is the dominating approach in topic modeling.

Bayesian learning is very powerful but too general theoretical framework, for which topic modeling is one of example applications. Bayesian inference is elegant when conjugate priors are used. However, the Dirichlet conjugate prior is not always a better choice from the natural language modeling point of view. In particular, it conflicts with natural assumptions of sparsity. Better motivated non-conjugate priors require a laborious mathematical work and lead to intricate learning algorithms. The development of combined and multi-objective topic models also remains a challenging task in Bayesian approach. An evolutionary approach to multi-objective Bayesian topic modeling has been proposed in [13], but it seems to be computationally infeasible for large text collections. Until now, there was no freely available software to combine topic models.

Additive Regularization of Topic Models (ARTM) is a semi-probabilistic approach based on classical (non-Bayesian) regularization [26]. In ARTM a topic model is learned by maximizing a weighted sum of the log-likelihood and additional regularization criteria. These criteria, or

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

KDD KDD 2015

Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$15.00.

penalty terms, are not required to be log-priors or even to have a probabilistic sense. The optimization problem is solved by a general regularized expectation-maximization (EM) algorithm, which can be easily applied to any combination of regularization criteria. The non-Bayesian regularization provides a much simpler inference for many topic models previously studied only in the Bayesian setting [28, 27]. In particular, the LDA model can be alternatively understood as a smoothing regularizer that minimizes Kullback–Leibler divergence of each topic distribution with a fixed multinomial distribution. The maximization of the same Kullback–Leibler divergence naturally leads to a sparsing regularizer [28]. This possibility is difficult to see from the Bayesian perspective, thereby all Bayesian approaches to sparsing are much more complicated [23, 29, 14, 10, 5].

ARTM makes topic models easier to design, to explain, to infer, and to combine, thereby reducing barriers to entry into topic modeling research field.

In this paper we develop the ARTM approach in three directions: the multiple modalities modeling, the online algorithm, and the parallel implementation.

Multimodal data has become increasingly important in many application areas. Large collections of data coming from the web consist of heterogeneous linked data. Typically, texts are accompanied by images, audio or video clips, usage data, metadata containing authors, links, date-time stamps, etc. In these cases documents are considered as multimodal containers, for which words are the elements of only one of multiple modalities. All modalities are useful for determining more relevant topics, and, vice-versa, topics are useful for crossmodal retrieval or making predictions when data of some modalities are missing. We introduce the regularized multimodal topic model with an arbitrary number of modalities and then generalize the regularized EM-algorithm for this case.

Online algorithms have proven to be the fastest for the large document collections, including those arriving in a stream. Online algorithms are now available for PLSA [1], LDA variational inference [11], LDA stochastic inference [16], and some other topic models. We show that the online algorithm is not necessarily associated with a particular type of model, nor a particular type of inference, but only with a certain reorganization of steps in the EM-like iterative process. Our online algorithm remains the same for any combination of regularizers and any number of modalities.

The parallel implementation of topic model learning algorithm as well as the distributed storage of data are very important technical issues for the analysis of large document collections. Our parallel distributed implementation of the additively regularized multimodal online topic models is freely available as BigARTM open source project <http://bigartm.org>. BigARTM source code is released under the New BSD License, which permits free commercial and non-commercial usage. The core of the library is written in C++ and is exposed via two equally rich APIs for C++ and Python. The library is cross-platform and can be built for Linux, Windows and OS X in both 32 and 64 bit configuration.

The rest of the paper is organized as follows.

In section 2 we introduce notation and definitions of topic modeling and additive regularization.

In section 3 we introduce a multimodal topic modeling for

documents with additional discrete metadata.

In section 4 we generalize the fast online algorithm [11] to additively regularized multimodal topic models.

In section 5 we describe parallel architecture and implementation details of the BigARTM library.

In section 6 we report results of our experiments on large datasets. BigARTM performs better than Vowpal Wabbit LDA and Gensim libraries in terms of perplexity and runtime on Wikipedia corpus. Comparing to the other libraries BigARTM offers several additional features, such as regularization and multimodality.

In section 7 we discuss advantages, limitations and open problems of BigARTM.

2. ADDITIVE REGULARIZATION FOR PROBABILISTIC TOPIC MODELS

Let D denote a finite set (collection) of texts and W denote a finite set (vocabulary) of all terms from these texts. Each term can represent a single word or a key phrase. Each document $d \in D$ is a sequence of terms from the vocabulary W . Assume that each term occurrence in each document refers to some latent topic from a finite set of topics T . Text collection is considered to be a sample of triples (w_i, d_i, t_i) , $i = 1, \dots, n$, drawn independently from a discrete distribution $p(w, d, t)$ over the finite probability space $W \times D \times T$. Terms w_i and documents d_i are observable variables, while topics t_i are latent variables.

The topic model of Probabilistic Latent Semantic Analysis, PLSA [12] explains the terms probabilities $p(w | d)$ in each document $d \in D$ by a mixture of term probabilities for topics and topic probabilities for documents:

$$p(w | d) = \sum_{t \in T} p(w | t) p(t | d) = \sum_{t \in T} \phi_{wt} \theta_{td}, \quad w \in W.$$

This representation follows immediately from the law of total probability and the assumption of conditional independence $p(w | t) = p(w | d, t)$, which means that each topic generates terms regardless of the document.

The parameters $\theta_{td} = p(t | d)$ and $\phi_{wt} = p(w | t)$ form matrices $\Theta = (\theta_{td})_{T \times D}$ and $\Phi = (\phi_{wt})_{W \times T}$. These matrices are *stochastic*, that is, their vector-columns represent discrete distributions. The number of topics $|T|$ is usually much smaller than $|D|$ and $|W|$.

To learn parameters Φ, Θ from the collection we maximize the log-likelihood:

$$\mathcal{L}(\Phi, \Theta) = \sum_{d \in D} \sum_{w \in W} n_{dw} \ln p(w | d) \rightarrow \max_{\Phi, \Theta},$$

where n_{dw} is the number of occurrences of the term $w \in W$ in the document d .

Following the ARTM approach, we introduce r additional criteria $R_i(\Phi, \Theta)$, $i = 1, \dots, r$, called *regularizers*. We would like to maximize them separately, but the maximization of their linear combination with nonnegative *regularization coefficients* ρ_i is technically more convenient:

$$R(\Phi, \Theta) = \sum_{i=1}^r \rho_i R_i(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}.$$

Then we add a regularization penalty term $R(\Phi, \Theta)$ to the log-likelihood and solve a constrained multicriteria op-

timization problem via scalarization:

$$\mathcal{L}(\Phi, \Theta) + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}; \quad (1)$$

$$\sum_{w \in W} \phi_{wt} = 1, \phi_{wt} \geq 0; \quad \sum_{t \in T} \theta_{td} = 1, \theta_{td} \geq 0. \quad (2)$$

The local maximum (Φ, Θ) of the problem (6), (7) satisfies the following system of equations with auxiliary variables $p_{tdw} = p(t|d, w)$, which follows from Karush–Kuhn–Tucker conditions [28]:

$$p_{tdw} = \text{norm}_{t \in T}(\phi_{wt}\theta_{td}); \quad (3)$$

$$n_{wt} = \sum_{d \in D} n_{dw} p_{tdw};$$

$$n_{td} = \sum_{w \in d} n_{dw} p_{tdw};$$

$$\phi_{wt} = \text{norm}_{w \in W} \left(n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right); \quad (4)$$

$$\theta_{td} = \text{norm}_{t \in T} \left(n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right); \quad (5)$$

where operator $\text{norm}_{t \in T} x_t = \frac{\max\{x_t, 0\}}{\sum_{s \in T} \max\{x_s, 0\}}$ transforms a vector $(x_t)_{t \in T}$ to a discrete distribution.

The system of equations (3)–(5) can be solved by various numerical methods. In particular, the simple-iteration method is equivalent to the EM algorithm, which is typically used in practice.

Many Bayesian topic models can be considered as special cases of ARTM with different regularizers [28]. For example, PLSA [12] corresponds to the absence of regularization, $R = 0$. LDA [4] corresponds to the smoothing regularizer, which minimizes the KL-divergences $KL(\alpha||\theta_d)$ and $KL(\beta||\phi_t)$ for fixed distributions β, α . The choice of these distributions as uniform corresponds to the use of symmetric Dirichlet priors in Bayesian approach.

Due to the additivity ARTM can build topic models for various applications simply by choosing a suitable combination of predefined regularizers from a user extendable library.

In experiments described below we use a combination of five regularizers that helps to improve the interpretability of topics. First, we split the set of topics T into two subsets, S and B . Domain-specific topics $t \in S$ contains terms of particular domain areas. It is desirable for them to be sparse and weakly correlated. To make them sparse we maximize the KL-divergences $KL(\alpha||\theta_d)$ and $KL(\beta||\phi_t)$ as regularizers. Background topics $t \in B$ contains common lexis words. They must be smooth because background words occur in many documents. To make them smooth we minimize the KL-divergences $KL(\alpha||\theta_d)$ and $KL(\beta||\phi_t)$. Also we introduce the covariance regularizer to make all topics weakly correlated as columns of Φ matrix. The additive combination of regularizers is summarized as follows:

$$\begin{aligned} R(\Phi, \Theta) = & -\beta_0 \sum_{t \in S} \sum_{w \in W} \beta_w \ln \phi_{wt} - \alpha_0 \sum_{d \in D} \sum_{t \in S} \alpha_t \ln \theta_{td} \\ & + \beta_1 \sum_{t \in B} \sum_{w \in W} \beta_w \ln \phi_{wt} + \alpha_1 \sum_{d \in D} \sum_{t \in B} \alpha_t \ln \theta_{td} \\ & - \gamma \sum_{t \in T} \sum_{s \in T \setminus t} \sum_{w \in W} \phi_{wt} \phi_{ws}, \end{aligned}$$

where $\beta_0, \alpha_0, \beta_1, \alpha_1, \gamma$ are regularization coefficients.

3. MULTIMODAL REGULARIZED TOPIC MODEL

Now assume that a document can contain not only words, but also terms of other modalities. Each modality is defined by a finite set (vocabulary) of terms W^m , $m = 1, \dots, M$.

Examples of not-word modalities are: authors, class or category labels, date-time stamps, references to/from other documents/authors, named entities mentioned in texts, objects found in the images associated with the documents, users that read or downloaded documents, advertising banners, etc.

As in the previous section, the collection is considered to be a sample of i.i.d. triples $(w_i, d_i, t_i) \sim p(w, d, t)$ drawn from the finite probability space $W \times D \times T$, but now $W = W^1 \sqcup \dots \sqcup W^m$ is a disjoint union of the vocabularies across all modalities.

Following the idea of Correspondence LDA [3] and Dependency LDA [22] we introduce a topic model $p(w|d)$ for each modality W^m , $m = 1, \dots, M$:

$$p(w|d) = \sum_{t \in T} p(w|t) p(t|d) = \sum_{t \in T} \phi_{wt} \theta_{td}, \quad w \in W^m.$$

Stochastic matrices $\Phi^m = (\phi_{wt})_{W^m \times T}$ of term probabilities for the topics, if stacked vertically, form a $W \times T$ -matrix Φ .

To learn parameters Φ^m, Θ from the multimodal collection we maximize the log-likelihood for each m -th modality:

$$\mathcal{L}_m(\Phi^m, \Theta) = \sum_{d \in D} \sum_{w \in W^m} n_{dw} \ln p(w|d) \rightarrow \max_{\Phi^m, \Theta},$$

where n_{dw} is the number of occurrences of the term $w \in W^m$ in the document d . Note that topic distributions of documents Θ are common for all modalities.

In ARTM we add a regularization penalty term $R(\Phi, \Theta)$ to the log-likelihood and solve a constrained multicriteria optimization problem:

$$\sum_{m=1}^M \tau_m \mathcal{L}_m(\Phi^m, \Theta) + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}; \quad (6)$$

$$\sum_{w \in W^m} \phi_{wt} = 1, \phi_{wt} \geq 0; \quad \sum_{t \in T} \theta_{td} = 1, \theta_{td} \geq 0; \quad (7)$$

where *regularization coefficients* τ_m are used to balance the importance of different modalities.

The local maximum (Φ, Θ) of the problem (6), (7) satisfies the following system of equations with auxiliary variables $p_{tdw} = p(t|d, w)$:

$$p_{tdw} = \text{norm}_{t \in T}(\phi_{wt}\theta_{td}); \quad (8)$$

$$n_{wt} = \sum_{d \in D} n_{dw} p_{tdw};$$

$$n_{td} = \sum_{w \in d} \tau_{m(w)} n_{dw} p_{tdw};$$

$$\phi_{wt} = \text{norm}_{w \in W^m} \left(n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right); \quad (9)$$

$$\theta_{td} = \text{norm}_{t \in T} \left(n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right); \quad (10)$$

where $m(w)$ is the modality of the term w , $w \in W^{m(w)}$.

The system of equations (8)–(10) follows from Karush–Kuhn–Tucker conditions (see Appendix A for the proof).

For single modality ($M = 1$) it gives the regularized EM algorithm described in the previous section.

Many previous topic models for labeled documents can be considered as special cases of multimodal ARTM. Most of them are based on LDA model and use Dirichlet priors, which correspond to smoothing regularization. From ARTM perspective, there is little reason to always use only the smoothing regularizer.

Following topic models exactly correspond to the multimodal ARTM, up to the modality sense. A topic model of document content and hypertext connectivity [6] has the modality of documents to which a given document has a hyperlink. The Conditionally Independent LDA, CI-LDA [19] has the modality of named entities mentioned in a given document. The Tag-LDA [24] has the modality of tags as special kind of words. The LDA-JS and LDA-post [9] has the modality of publications cited in a given document; an additional regularizer takes into account that cited documents are likely to share similar topics. Both models are designed to estimate the strength of influence of cited publications. The Dependency LDA [22] has the modality of document categories or class labels. The MultiLingual LDA, ML-LDA [20] and the PolyLingual Topic Model, PLTM [17] have L modalities for L different languages; parallel documents always share one identical topic distribution. The BiLingual LDA, BiLDA [8] is also a multilanguage topic model, but the number of modalities is restricted by two.

4. ONLINE TOPIC MODELING

Like Online LDA [11] and Online PLSA [1] we split the collection D into batches D_b , $b = 1, \dots, B$, and organize EM iterations so that each document vector θ_d is iterated until convergence at a constant matrix Φ , see Algorithm 1 and 2. Matrix Φ is updated rarely, after all documents from the batch are processed. For a large collection matrix Φ often stabilizes after small initial part of the collection. Therefore a single pass through the collection might be sufficient to learn a topic model. The second pass may be needed for the initial part of the collection.

Algorithm 1 does not specify how often to synchronize Φ matrix at steps 5–8. It can be done after every batch or less frequently (for instance if $\frac{\partial R}{\partial \phi_{wt}}$ takes long time to evaluate). This flexibility is especially important for concurrent implementation of the algorithm, where multiple batches are processed in parallel. In this case synchronization can be triggered when a fixed number of documents had been processed since the last synchronization.

The online reorganization of the EM iterations is not necessarily associated with Bayesian inference used in [11]. Different topic models, from PLSA to multimodal and regularized models, can be learned by the above online EM algorithm.

5. BIGARTM ARCHITECTURE

The main goal for BigARTM architecture is to ensure a constant memory usage regardless of the collection size. For this reason each D_b batch is stored on disk in a separate file, and only a limited number of batches is loaded into the main memory at any given time. The entire Θ matrix is also never stored in the memory. As a result, the memory usage stays constant regardless of the size of the collection.

Algorithm 1: Online EM-algorithm for multimodal ARTM

Input: collection D_b , discounting factor $\rho \in (0, 1]$;
Output: matrix Φ ;

- 1 initialize ϕ_{wt} for all $w \in W$ and $t \in T$;
- 2 $n_{wt} := 0$, $\tilde{n}_{wt} := 0$ for all $w \in W$ and $t \in T$;
- 3 **for all** batches D_b , $b = 1, \dots, B$
- 4 $(\tilde{n}_{wt}) := (\tilde{n}_{wt}) + \text{ProcessBatch}(D_b, \phi_{wt})$;
- 5 **if** (*synchronize*) **then**
- 6 $n_{wt} := \rho n_{wt} + \tilde{n}_{wt}$ for all $w \in W$ and $t \in T$;
- 7 $\phi_{wt} := \text{norm}_{w \in W^m} (n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}})$ for all $w \in W^m$,
- 8 $m = 1, \dots, M$ and $t \in T$;
- 8 $\tilde{n}_{wt} := 0$ for all $w \in W$ and $t \in T$;

Algorithm 2: ProcessBatch(D_b, ϕ_{wt})

Input: batch D_b , matrix ϕ_{wt} ;
Output: matrix (\tilde{n}_{wt}) ;

- 1 $\tilde{n}_{wt} := 0$ for all $w \in W$ and $t \in T$;
- 2 **for all** $d \in D_b$
- 3 initialize $\theta_{td} := \frac{1}{|T|}$ for all $t \in T$;
- 4 **repeat**
- 5 $p_{tdw} := \text{norm}_{t \in T} (\phi_{wt} \theta_{td})$ for all $t \in T$;
- 6 $n_{td} := \sum_{w \in d} \tau_m(w) n_{dw} p_{tdw}$ for all $t \in T$;
- 7 $\theta_{td} := \text{norm}_{t \in T} (n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}})$ for all $t \in T$;
- 8 **until** θ_d converges;
- 9 increment \tilde{n}_{wt} by $n_{dw} p_{tdw}$ for all $w \in d$ and $t \in T$;

Concurrency.

A general rule of concurrency design is to express parallelism at the highest possible level. For this reason BigARTM implements a concurrent processing of the batches and keeps a single-threaded code for the $\text{ProcessBatch}(D_b, \phi_{wt})$ routine.

To split collection into batches and process them concurrently is a common approach, introduced in AD-LDA algorithm [18], and then further developed in PLDA [30] and PLDA+ [15] algorithms. These algorithms require all concurrent workers to become idle before an update of the Φ matrix. Such synchronization step adds a large overhead in the online algorithm where Φ matrix is updated multiple times on each iteration. An alternative architecture without the synchronization step is described in [25], however it mostly targets a distributed cluster environment. In our work we develop an efficient single-node architecture where all workers benefit from the shared memory space.

To run multiple ProcessBatch in parallel the inputs and outputs of this routine are stored in two separate in-memory queues, locked for push and pop operations with spin locks. This approach does not add any noticeable synchronization overhead because both queues only store smart pointers to the actual data objects, so push and pop operations does not involve copying or relocating big objects in the memory.

Smart pointers are also essential for lifecycle of the Φ matrix. This matrix is *read* by all processors threads, and can be *written* at any time by the merger thread. To update Φ without pausing all processor threads we keep two copies —

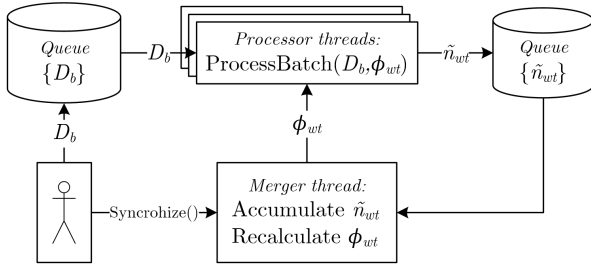


Figure 1: Diagram of key BigARTM components

an *active* Φ and a *background* Φ matrices. The active matrix is read-only, and is used by the processor threads. The background matrix is being built in a background by the merger thread at steps 6 and 7 of Algorithm 1, and once it is ready merger thread marks it as active. Before processing a new batch the processor thread gets the current active matrix from the merger thread. This object is passed via shared smart pointer to ensure that processor thread can keep ownership of its Φ matrix until the batch is fully processed. As a result, all processor threads keep running concurrently with the update of Φ matrix.

Note that all processor threads share the same Φ matrix, which means that memory usage stays at constant level regardless of how many cores are used for computation. Using memory for two copies of the Φ matrix in our opinion gives a reasonable usage balance between memory and CPU resources. An alternative solution with only one Φ matrix is also possible, but it would require a heavy usage of atomic CPU instructions. Such operations are very efficient, but still come at a considerable synchronization cost¹, and using them for all reads and writes of the Φ matrix would cause a significant performance degradation for merger and processor threads. Besides, an arbitrary overlap between reads and writes of the Φ matrix eliminates any possibility of producing a deterministic result. The design with two copies of the Φ matrix gives much more control over this and in certain cases allows BigARTM to behave in a fully deterministic way.

The design with two Φ matrices only supports a single merger thread, and we believe it should handle all \tilde{n}_{wt} updates coming from many threads. This is a reasonable assumption because merging at step 6 takes only about $O(|W| \cdot |T|)$ operations to execute, while `ProcessBatch` takes $O(n|T|I)$ operations, where n is the number of non-zero entries in the batch, I is the average number of inner iterations in `ProcessBatch` routine. The ratio $n/|W|$ is typically from 100 to 1000 (based on datasets in UCI Bag-Of-Words repository), and I is 10...20, so the ratio safely exceeds the expected number of cores (up to 32 physical CPU cores in modern workstations, and even 60 cores of the Intel Xeon Phi co-processors).

Data layout.

BigARTM uses dense single-precision matrices to represent Φ and Θ . Together with the Φ matrix we store a global dictionary of all terms $w \in W$. This dictionary is implemented as `std::unordered_map` that maps a string rep-

resentation of $w \in W$ into its integer index in the Φ matrix. This dictionary can be extended automatically as more and more batches came through the system. To achieve this each batch D_b contains a local dictionary W_b , listing all terms that occur in the batch. The n_{dw} elements of the batch are stored as a sparse CSR matrix (Compressed Sparse Row format), where each row correspond to a document $d \in D_b$, and terms w run over a local batch dictionary W_b .

For performance reasons Φ matrix is stored in column-major order, and Θ in row-major order. This layout ensures that $\sum_t \phi_{wt} \theta_{td}$ sum runs on contiguous memory blocks. In both matrices all values smaller than 10^{-16} are always replaced with zero to avoid performance issues with denormalized numbers².

Programming interface.

All functionality of BigARTM is expressed in a set of extern C methods. To input and output complex data structures the API uses Google Protocol Buffers³. This approach makes it easy to integrate BigARTM into any research or production environment, as almost every modern language has an implementation of Google Protocol Buffers and a way of calling extern C code (ctypes module for Python, loadlibrary for Matlab, Plnvoke for C#, etc).

On top of the extern C API BigARTM already has convenient wrappers in C++ and Python. We are also planning to implement a Java wrapper in the near future. In addition to the APIs the library also has a simple CLI interface.

BigARTM has built-in libraries of regularizers and quality measures that can be extended in current implementation only through project recompilation.

Basic tools.

A careful selection of the programming tools is important for any software project. This is especially true for BigARTM as its code is written in C++, a language that by itself offers less functionality comparing to Python, .NET Framework or Java. To mitigate this we use various parts of the Boost C++ Libraries, Google Protocol Buffers for data serialization, ZeroMQ library for network communication, and several other libraries.

BigARTM uses CMake as a cross-platform build system, and it successfully builds on Windows, Linux and OS X in 32 and 64 bit configurations. Building the library require a recent C++ compiler with C++11 support (GNU GCC 4.6.3, clang 3.4 or Visual Studio 2012 or newer), and Boost Libraries 1.46.1 or newer. All the other third-parties are included in BigARTM repository.

We also use free online services to store source code (<https://github.com>), to host online documentation (<https://readthedocs.org>) and to run automated continuous integration builds (<http://travis-ci.org>).

6. EXPERIMENTS

In this section we evaluate the runtime performance and the algorithmic quality of BigARTM against two popular software packages — Gensim [21] and Vowpal Wabbit⁴. We also demonstrate some of the unique BigARTM features,

²http://en.wikipedia.org/wiki/Denormal_number#Performance_issues

³<http://code.google.com/p/protobuf/>

⁴https://github.com/JohnLangford/vowpal_wabbit/

¹<http://stackoverflow.com/questions/2538070/atomic-operation-cost>

such as combining regularizers and multi-language topic modeling via multimodality, which are not available in the other software packages.

All three libraries (VW.LDA, Gensim and BigARTM) work out-of-core, e.g. they are designed to process data that is too large to fit into a computer’s main memory at one time. This allowed us to benchmark on a fairly large collection — 3.7 million articles from the English Wikipedia⁵. The conversion to bag-of-words was done with `gensim.make_wikicorpus` script⁶, which excludes all non-article pages (such as category, file, template, user pages, etc), and also pages that contain less than 50 words. The dictionary is formed by all words that occur in at least 20 documents, but no more than in 10% documents in the collection. The resulting dictionary was capped at $|W| = 100\,000$ most frequent words.

Both Gensim and VW.LDA represents the resulting topic model as Dirichlet distribution over Φ and Θ matrices: $\theta_d \sim \text{Dir}(\gamma_d)$ and $\phi_t \sim \text{Dir}(\lambda_t)$. On contrary, BigARTM outputs a non-probabilistic matrices Φ and Θ . To compare the perplexity we take the mean or the mode of the posterior distributions:

$$\begin{aligned}\phi_{wt}^{\text{mean}} &= \text{norm } \lambda_{wt}, & \theta_{td}^{\text{mean}} &= \text{norm } \gamma_{td}; \\ \phi_{wt}^{\text{mode}} &= \text{norm}(\lambda_{wt} - 1), & \theta_{td}^{\text{mode}} &= \text{norm}(\gamma_{td} - 1).\end{aligned}$$

The perplexity measure is defined as

$$\mathcal{P}(D, p) = \exp\left(-\frac{1}{n} \sum_{d \in D} \sum_{w \in d} n_{dw} \ln p(w|d)\right). \quad (11)$$

Comparison to existing software packages.

The *Vowpal Wabbit* (VW) is a library of online algorithms that cover a wide range of machine learning problems. For topic modeling VW has the VW.LDA algorithm, based on the Online Variational Bayes LDA [11]. VW.LDA is neither multi-core nor distributed, but an effective single-threaded implementation in C++ made it one of the fastest tools for topic modeling.

The *Gensim* library specifically targets the area of topic modeling and matrix factorization. It has two LDA implementations — `LdaModel` and `LdaMulticore`, both based on the same algorithm as VW.LDA (Online Variational Bayes LDA [11]). Gensim is entirely written in Python. Its high performance is achieved through the usage of NumPy library, built over low-level BLAS libraries (such as Intel MKL, ATLAS, or OpenBLAS). In `LdaModel` all batches are processed sequentially, and the concurrency happens entirely within NumPy. In `LdaMulticore` the workflow is similar to BigARTM — several batches are processed concurrently, and there is a single aggregation thread that asynchronously merges the results.

Each run in our experiment performs one pass over the Wikipedia corpus and produces a model with $|T| = 100$ topics. The runtime is reported for an Intel-based CPU with 16 physical cores with hyper-threading. The collection was split into batches with 10000 documents each (`chunksize` in Gensim, `minibatch` in VW.LDA). The update rule in online algorithm used $\rho = (b + \tau_0)^{-0.5}$, where b is the number

⁵<http://dumps.wikimedia.org/enwiki/20141208/>

⁶<https://github.com/piskvorky/gensim/tree/develop/gensim/scripts/>

Table 1: The comparison of BigARTM with VW.LDA and Gensim. *Train time* is the time for model training, *inference* is the time for calculation of θ_d of 100 000 held-out documents, *perplexity* is calculated according to (11) on held-out documents.

library	procs	train time	inference time	perplexity	
				mode	mean
BigARTM	1	35 min	72 sec	4000	
LdaModel	1	369 min	395 sec	4213	4161
VW.LDA	1	73 min	120 sec	4061	4108
BigARTM	4	9 min	20 sec	4061	
LdaMulticore	4	60 min	222 sec	4055	4111
BigARTM	8	4.5 min	14 sec	4304	
LdaMulticore	8	57 min	224 sec	4379	4455

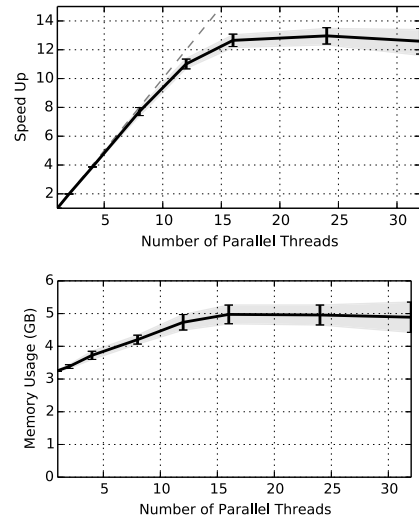


Figure 2: Running BigARTM in parallel: speed up (upper chart) and memory usage (lower chart)

of batches processed so far, and τ_0 is a constant offset parameter introduced in [11], in our experiment $\tau_0 = 64$. Updates were performed after each batch in non-parallel runs, and after P batches when running in P threads. LDA priors were fixed as $\alpha = 0.1$, $\beta = 0.1$, so that $\theta_d \sim \text{Dir}(\alpha)$, $\phi_t \sim \text{Dir}(\beta)$.

Table 1 compares the performance of VW.LDA, Gensim `LdaModel`, Gensim `LdaMulticore`, and BigARTM. Fig. 2 shows BigARTM speedup and memory consumption depending on the number of CPU threads for Amazon AWS c3.8xlarge with 32 virtual cores, Gensim 0.10.3 under Python 2.7.

Experiments with combination of regularizers.

BigARTM has a built-in library of regularizers, which can be used in any combination. In the following experiment we combine tree regularizers described at the end of section 2: sparsing of ϕ_t distributions, sparsing of θ_d distributions, and pairwise decorrelation of ϕ_t distributions. This combination helps to improve several quality measures without significant loss of perplexity, according to experiments on the offline implementation of ARTM [28]. The goal of our experiment is

Table 2: Comparison of LDA and ARTM models. Quality measures: \mathcal{P}_{10k} , \mathcal{P}_{100k} — hold-out perplexity on 10K and 100K documents sets, \mathcal{S}_Φ , \mathcal{S}_Θ — sparsity of Φ and Θ matrices (in %), \mathcal{K}_s , \mathcal{K}_p , \mathcal{K}_c — average topic kernel size, purity and contrast respectively.

Model	\mathcal{P}_{10k}	\mathcal{P}_{100k}	\mathcal{S}_Φ	\mathcal{S}_Θ	\mathcal{K}_s	\mathcal{K}_p	\mathcal{K}_c
LDA	3436	3801	0.0	0.0	873	0.533	0.507
ARTM	3577	3947	96.3	80.9	1079	0.785	0.731

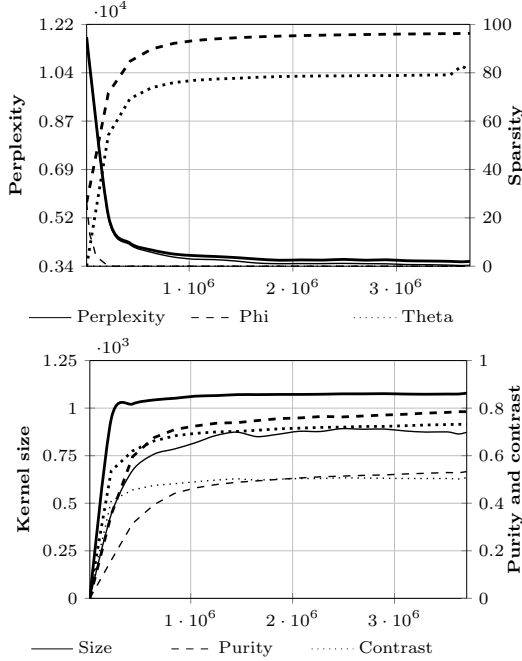


Figure 3: Comparison of LDA (thin) and ARTM (bold) models. The number of processed documents is shown along the X axis.

to show that this remains true for the online implementation in BigARTM. We use the following built-in quality measures: the hold-out perplexity, the sparsity of Φ and Θ matrices, and the characteristics of topic lexical kernels (size, purity, and contrast) averaged across all topics.

Table 2 compares the results of additive combination of regularizers (ARTM) and the usual LDA model. Figure 3 presents quality measures as functions of the number of processed documents. The first chart shows perplexity and sparsity of Φ , Θ matrices, and the second chart shows average lexical kernel measures.

Experiments on multi-language Wikipedia.

To show how BigARTM works with multimodal datasets we prepared a text corpus containing all English and Russian Wikipedia articles with mutual interlanguage links. We represent each linked pair of articles as a single multi-language document with two modalities, one modality for each language. That is how our multi-language collection acts as a multimodal document collection.

The dump of Russian articles⁷ had been processed follow-

ing the same technique as we previously used in experiments on English Wikipedia. Russian words were lemmatized with Yandex MyStem 3.0⁸. To further reduce the dictionary we only keep words that appear in no less than 20 documents, but no more than in 10% of documents in the collection. The resulting collection contains 216175 pairs of Russian–English articles, with combined dictionary of 196749 words (43% Russian, 57% English words).

We build multi-language model with 400 topics. They cover a wide range of themes such as science, architecture, history, culture, technologies, army, different countries. All 400 topics were reviewed by an independent assessor, and he successfully interpreted all except four topics.

7. CONCLUSIONS

BigARTM in an open source project for parallel online topic modeling of large text collections. It provides a high flexibility for various applications due to multimodality and additive combinations of regularizers. BigARTM architecture has a rich potential. Current components can be reused in a distributed solution that runs on cluster. Further improvement of single-node can be achieved by offloading batch processing into GPU.

Acknowledgements.

The work was supported by the Russian Foundation for Basic Research grants 14-07-00847, 14-07-00908, 14-07-31176 and by Skolkovo Institute of Science and Technology (project 081-R).

8. REFERENCES

- [1] N. Bassiou and C. Kotropoulos. Online pls: Batch updating techniques including out-of-vocabulary words. *Neural Networks and Learning Systems, IEEE Transactions on*, 25(11):1953–1966, Nov 2014.
- [2] D. M. Blei. Probabilistic topic models. *Communications of the ACM*, 55(4):77–84, 2012.
- [3] D. M. Blei and M. I. Jordan. Modeling annotated data. In *Proceedings of the 26th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 127–134, New York, NY, USA, 2003. ACM.
- [4] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022, 2003.
- [5] J.-T. Chien and Y.-L. Chang. Bayesian sparse topic model. *Journal of Signal Processing Systems*, 74:375–389, 2013.
- [6] D. A. Cohn and T. Hofmann. The missing link - a probabilistic model of document content and hypertext connectivity. In *NIPS*, pages 430–436, 2000.
- [7] A. Daud, J. Li, L. Zhou, and F. Muhammad. Knowledge discovery through directed probabilistic topic models: a survey. *Frontiers of Computer Science in China*, 4(2):280–301, 2010.
- [8] W. De Smet and M.-F. Moens. Cross-language linking of news stories on the web using interlingual topic modelling. In *Proceedings of the 2Nd ACM Workshop*

⁷<http://dumps.wikimedia.org/ruwiki/20141203/>

⁸<https://tech.yandex.ru/mystem/>

- on Social Web Search and Mining, SWSM '09, pages 57–64, New York, NY, USA, 2009. ACM.
- [9] L. Dietz, S. Bickel, and T. Scheffer. Unsupervised prediction of citation influences. In *Proceedings of the 24th international conference on Machine learning, ICML '07*, pages 233–240, New York, NY, USA, 2007. ACM.
 - [10] J. Eisenstein, A. Ahmed, and E. P. Xing. Sparse additive generative models of text. In *ICML'11*, pages 1041–1048, 2011.
 - [11] M. D. Hoffman, D. M. Blei, and F. R. Bach. Online learning for latent dirichlet allocation. In *NIPS*, pages 856–864. Curran Associates, Inc., 2010.
 - [12] T. Hofmann. Probabilistic latent semantic indexing. In *Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval*, pages 50–57, New York, NY, USA, 1999. ACM.
 - [13] O. Khalifa, D. Corne, M. Chantler, and F. Halley. Multi-objective topic modelling. In *7th International Conference Evolutionary Multi-Criterion Optimization (EMO 2013)*, pages 51–65. Springer LNCS, 2013.
 - [14] M. O. Larsson and J. Ugander. A concave regularization technique for sparse mixture models. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K. Weinberger, editors, *Advances in Neural Information Processing Systems 24*, pages 1890–1898, 2011.
 - [15] Z. Liu, Y. Zhang, E. Y. Chang, and M. Sun. PLDA+: parallel latent Dirichlet allocation with data placement and pipeline processing. *ACM Trans. Intell. Syst. Technol.*, 2(3):26:1–26:18, May 2011.
 - [16] D. Mimno, M. Hoffman, and D. Blei. Sparse stochastic inference for latent dirichlet allocation. In J. Langford and J. Pineau, editors, *Proceedings of the 29th International Conference on Machine Learning (ICML-12)*, pages 1599–1606, New York, NY, USA, July 2012. Omnipress.
 - [17] D. Mimno, H. M. Wallach, J. Naradowsky, D. A. Smith, and A. McCallum. Polylingual topic models. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 2 - Volume 2*, EMNLP '09, pages 880–889, Stroudsburg, PA, USA, 2009. Association for Computational Linguistics.
 - [18] D. Newman, A. Asuncion, P. Smyth, and M. Welling. Distributed algorithms for topic models. *J. Mach. Learn. Res.*, 10:1801–1828, Dec. 2009.
 - [19] D. Newman, C. Chemudugunta, and P. Smyth. Statistical entity-topic models. In *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '06*, pages 680–686, New York, NY, USA, 2006. ACM.
 - [20] X. Ni, J.-T. Sun, J. Hu, and Z. Chen. Mining multilingual topics from wikipedia. In *Proceedings of the 18th International Conference on World Wide Web, WWW '09*, pages 1155–1156, New York, NY, USA, 2009. ACM.
 - [21] R. Řehůřek and P. Sojka. Software framework for topic modelling with large corpora. In *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, pages 45–50, Valletta, Malta, May 2010. ELRA.
 - [22] T. N. Rubin, A. Chambers, P. Smyth, and M. Steyvers. Statistical topic models for multi-label document classification. *Machine Learning*, 88(1-2):157–208, 2012.
 - [23] M. Shashanka, B. Raj, and P. Smaragdis. Sparse overcomplete latent variable decomposition of counts data. In J. C. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems, NIPS-2007*, pages 1313–1320. MIT Press, Cambridge, MA, 2008.
 - [24] X. Si and M. Sun. Tag-lda for scalable real-time tag recommendation. *Journal of Information & Computational Science*, 6:23–31, 2009.
 - [25] A. Smola and S. Narayanamurthy. An architecture for parallel topic models. *Proc. VLDB Endow.*, 3(1-2):703–710, Sept. 2010.
 - [26] K. V. Vorontsov. Additive regularization for topic models of text collections. *Doklady Mathematics*, 89(3):301–304, 2014.
 - [27] K. V. Vorontsov and A. A. Potapenko. Additive regularization of topic models. *Machine Learning, Special Issue on Data Analysis and Intelligent Optimization*, 2014.
 - [28] K. V. Vorontsov and A. A. Potapenko. Tutorial on probabilistic topic modeling: Additive regularization for stochastic matrix factorization. In *AIST'2014, Analysis of Images, Social networks and Texts*, volume 436, pages 29–46. Springer International Publishing Switzerland, Communications in Computer and Information Science (CCIS), 2014.
 - [29] C. Wang and D. M. Blei. Decoupling sparsity and smoothness in the discrete hierarchical dirichlet process. In *NIPS*, pages 1982–1989. Curran Associates, Inc., 2009.
 - [30] Y. Wang, H. Bai, M. Stanton, W.-Y. Chen, and E. Y. Chang. PLDA: Parallel latent Dirichlet allocation for large-scale applications. In *Proceedings of the 5th International Conference on Algorithmic Aspects in Information and Management, AAIM '09*, pages 301–314, Berlin, Heidelberg, 2009. Springer-Verlag.

Appendix A

Consider the system of equations (8)–(10).

Topic t is called *regular* for the modality m if $n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} > 0$ for at least one term $w \in W^m$. If the reverse inequality holds for all $w \in W^m$ then topic t is called *irregular*.

Document d is called *regular* if $n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} > 0$ for at least one topic $t \in T$. If the reverse inequality holds for all $t \in T$ then document d is called *irregular*.

THEOREM 1. *If the function $R(\Phi, \Theta)$ is continuously differentiable and (Φ, Θ) is the local maximum of the problem (6), (7) then for any regular topic-modality pair (t, m) and any regular document d the system of equations (8)–(10) holds.*

Note 1. If a topic t is irregular then the t -th vector-column in matrix Φ^m equals zero and can not represent a discrete distribution. This means that topic t for the modality m must be excluded from the model. This mechanism is useful for irrelevant topics elimination and determining the number of topics.

Note 2. If a documents d is irregular then the d -th vector-column in matrix Θ equals zero and can not represent a discrete distribution. This means that document d must be excluded from the model. For example, a document may be too short or irrelevant to the given collection.

PROOF. For the local minimum Φ^m, Θ of the problem (6), (7) the Karush–Kuhn–Tucker (KKT) conditions can be written as follows:

$$\begin{aligned} \sum_d n_{dw} \frac{\theta_{td}}{p(w|d)} + \frac{\partial R}{\partial \phi_{wt}} &= \lambda_t - \lambda_{wt}; \\ \lambda_{wt} &\geq 0; \quad \lambda_{wt} \phi_{wt} = 0; \\ \sum_m \tau_m \sum_{w \in W^m} n_{dw} \frac{\phi_{wt}}{p(w|d)} + \frac{\partial R}{\partial \theta_{td}} &= \mu_d - \mu_{td}; \\ \mu_{td} &\geq 0; \quad \mu_{td} \theta_{td} = 0; \end{aligned}$$

where $\lambda_t, \lambda_{wt}, \mu_d, \mu_{td}$ are KKT multipliers for normalization and nonnegativity constraints.

Let us multiply both sides of the first equation by ϕ_{wt} , both sides of the second equation by θ_{td} , and reveal the auxiliary variable p_{tdw} from (8) in the left-hand side of both equations. Then we sum the right-hand side of the first equation over d , the right-hand side of the second equation over t :

$$\begin{aligned} \phi_{wt} \lambda_t &= \sum_d n_{dw} \frac{\phi_{wt} \theta_{td}}{p(w|d)} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} = n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}}; \\ \theta_{td} \mu_d &= \sum_m \tau_m \sum_{w \in W^m} n_{dw} \frac{\phi_{wt} \theta_{td}}{p(w|d)} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} = n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}}. \end{aligned}$$

An assumption that $\lambda_t \leq 0$ contradicts the regularity condition for the (t, m) pair. Then $\lambda_t > 0$. Either $\phi_{wt} = 0$ or both sides of the first equation are positive. Combining these two cases in one formula, we write:

$$\phi_{wt} \lambda_t = \max \left\{ n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}}, 0 \right\}. \quad (12)$$

Analogously, an assumption that $\mu_d \leq 0$ contradicts the regularity condition for the document d . Then $\mu_d > 0$. Either

$\theta_{td} = 0$ or both sides of the second equation are positive, consequently,

$$\theta_{td} \mu_d = \max \left\{ n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}}, 0 \right\}. \quad (13)$$

Let us sum both sides of the first equation over all $w \in W^m$, then both sides of the second equation over all $t \in T$:

$$\lambda_t = \sum_{w \in W^m} \max \left\{ n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}}, 0 \right\}; \quad (14)$$

$$\mu_d = \sum_{t \in T} \max \left\{ n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}}, 0 \right\}. \quad (15)$$

Finally, we obtain (9) by expressing ϕ_{wt} from (12) and (14).

Analogously, we obtain (10) by expressing θ_{td} from (13) and (15). \square